| REPORT DOCUMENTATION PAGE | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|

| 1. REPORT NUMBER | 2. GOVT ACCESSION NO. | 3. RECIPIENT'S CATALOG NUMBER |
|---|---|---|
| AFOSR-TR- 80-1092 AD-A100146 | | |

| 4. TITLE (and Subtitle) | 5. TYPE OF REPORT & PERIOD COVERED |
|---|---|
| MINIMAX RIDGE REGRESSION, | FINAL |
| | 6. PERFORMING ORG. REPORT NUMBER |

| 7. AUTHOR(s) | 8. CONTRACT OR GRANT NUMBER(s) |
|---|---|
| L. PEELE T.P. RYAN | F49620-79-C-0125 |

| 9. PERFORMING ORGANIZATION NAME AND ADDRESS | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS |
|---|---|
| Dept of Mathematics Sciences Old Dominion University Norfolk, VA | 61102F 2304/A5 |

| 11. CONTROLLING OFFICE NAME AND ADDRESS | 12. REPORT DATE |
|---|---|
| AFOSR/NM Bldg 410 Bolling AFB, DC 20332 | May 80 |
| | 13. NUMBER OF PAGES |
| | 22 |

| 14. MONITORING AGENCY NAME & ADDRESS(if different from Controlling Office) | 15. SECURITY CLASS. (of this report) |
|---|---|
| LEVEL | Unclassified |
| | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE |

**16. DISTRIBUTION STATEMENT (of this Report)**

Approved for Public Release; distribution unlimited

**17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)**

**18. SUPPLEMENTARY NOTES**

**19. KEY WORDS (Continue on reverse side if necessary and identify by block number)**

THIS DOCUMENT IS BEST QUALITY PRACTICABLE
THE COPY FURNISHED TO DDC CONTAINED A
SIGNIFICANT NUMBER OF PAGES WHICH DO NOT

**20. ABSTRACT (Continue on reverse side if necessary and identify by block number)**

This work examined minimax linear estimation in multiple linear regression. The application of minimax estimation to regression led to the development of ridge regression estimators with stochastic ridge parameters. These estimators were seen to be invariant under linear transformation: a property which has not been established for other ridge estimators. These minimax-motivated estimators were examined in several simulation studies. In particlar, flaws in other simulation studies of ridge estimators were depicted. (Continued)

DD FORM 1473 EDITION OF 1 NOV 65 IS OBSOLETE

# DISCLAIMER NOTICE

**THIS DOCUMENT IS BEST QUALITY PRACTICABLE. THE COPY FURNISHED TO DTIC CONTAINED A SIGNIFICANT NUMBER OF PAGES WHICH DO NOT REPRODUCE LEGIBLY.**

20. Cont.

Consequently, an improved simulation procedure was used. It was observed
from these studies that, contrary to published statements, a ridge estimator
can be considerably superior to the ordinary least squares estimator,
especially when high pairwise correlations exist among the regression variables.
Robustness considerations were used to suggest a requirement that a "good"
generalized ridge regression estimator should satisfy.

| Accession For | | |
|---|---|---|
| NTIS GRA&I | ☑ | |
| DTIC TAB | ☐ | |
| Unannounced | ☐ | |
| Justification | | |

| By | | |
|---|---|---|
| Distribution/ | | |
| Availability Codes | | |
| Dist | Avail and/or Special | |
| A | | |

MINIMAX RIDGE REGRESSION
(Final Report)

ABSTRACT

This work examined minimax linear estimation in multiple
linear regression. The application of minimax estimation to
regression led to the development of ridge regression estimators
with stochastic ridge parameters. These estimators were seen
to be invariant under linear transformation; a property which
has not been established for other ridge estimators. These
minimax-motivated estimators were examined in several simulation
studies. In particular, flaws in other simulation studies of
ridge estimators were depicted. Consequently, an improved
simulation procedure was used. It was observed from these
studies that, contrary to published statements, a ridge estimator
can be considerably superior to the ordinary least squares
estimator, especially when high pairwise correlations exist
among the regression variables. Robustness considerations were
used to suggest a requirement that a "good" generalized ridge
regression estimator should satisfy.

2. "Minimax Linear Regression Estimators with Application to Ridge Regression," by Lawrence Peele and Thomas P. Ryan, to appear in _Technometrics_.

3. "Comment" (on invited ridge regression paper) by Lawrence Peele and Thomas P. Ryan, _J. Amer. Statist. Assoc. 75_ 96-97.

4. "Most-robust Ridge Regression Estimators," by Lawrence Peele, submitted for publication.

# MINIMAX LINEAR REGRESSION ESTIMATORS
# WITH APPLICATION TO RIDGE REGRESSION[1]


by


Lawrence Peele and Thomas P. Ryan


AFOSR Technical Report No.1


May, 1980


Old Dominion University
Department of Mathematical Sciences
Norfolk, Virginia

## ABSTRACT

This paper considers minimax linear estimation of the parameters in a multiple linear regression model. Recent results are summarized, and some new results, including a transformation invariance property of minimax estimation, are given. These minimax estimators of the parameter vector can also be classified as ridge regression estimators with nonstochastic ridge parameters. Some ridge regression estimators with stochastic ridge parameters can be motivated by minimax estimation considerations. These minimax-motivated estimators are examined in several simulation studies and some observations are made based on these simulations and minimax theory.

## 1. INTRODUCTION

The usual multiple linear regression model is

$$Y = X\beta + \xi$$

where $\xi$ is a vector of uncorrelated random variables with mean zero and variance $\sigma^2$, and $X$ is a full rank $n \times q$ matrix with $q < n$. The usual procedure for estimating $\beta$ is the least squares method. It is well known that this method is equivalent to minimum variance unbiased linear (MVUL) estimation. Similarly, the usual method for estimating a given linear combination of the coefficients is MVUL. In recent years, many people have attempted to reduce the mean squared error (MSE) by allowing some bias in their estimators. One such biased estimation procedure is ridge regression, which was first studied by Hoerl and Kennard [7]. Ridge regression estimators have the form

$$\hat{\beta}^*(k) = [X'X + kI]^{-1}X'Y$$

where the ridge parameter $k$ is either nonstochastic (based on prior information) or stochastic. It follows from a result in [7] that for nonstochastic $k$ satisfying

$$k < \sigma^2/\beta'\beta \tag{1}$$

-2-

the ridge regression estimator $\hat{\beta}^*(k)$ outperforms the least squares estimator in terms of MSE. The minimax linear estimators discussed in the present paper are based on constraints related to (1). These estimators are shown (originally in [9]) to be ridge regression estimators, and both known and new properties of minimax linear regression estimators are given. Also, some simulation results are presented with an accompanying discussion.

## 2. MINIMAX LINEAR ESTIMATION

In general, minimax procedures "guard against the worst." It is observed in Section 3 that such caution is appropriate in highly ill-conditioned regression problems.

The following lemma is the basis for minimax linear estimation . The Cauchy-Schwarz inequality (See, for example, page 215 of [12].) when applied to $R^q$ results in the following lemma.

Lemma 1. Let $a \in R^q$. The maximum of $(a'z)^2$ among $z \in R^q$ satisfying $z'z \leq 1$ is $a'a$.

Let $k > 0$. Let $A$ be a positive definite (p.d.) matrix of order $q$, and consider the constraint

$$\beta' A \beta / \sigma^2 \leq k. \qquad (2)$$

Let $a \in R^q$ and suppose that we wish to estimate $a'\beta$. If $\hat{\gamma}$ satisfies

-3-

$$\text{...} \quad \max_{\gamma,\sigma \in (\Omega)} \quad E_{\beta,\sigma}[(\hat{\sigma}'Y - a'\beta)^2]$$

$$= \max_{\gamma,\sigma \in (\Omega)} \quad E_{\beta,\sigma}[(\tilde{\sigma}'Y - a'\beta)^2]$$

then $\hat{\sigma}'Y$ is the minimax linear estimator of $a'\beta$ based on (2).

The following result is presented in [9]. A brief proof is given as the method of proof seems useful.

Theorem 1. The minimax linear estimator $\hat{\sigma}'Y$ of $a'\beta$ based on (2) satisfies $\hat{\sigma}'Y = a'\hat{\beta}_{A,h}$ where

$$\hat{\beta}_{A,h} = [X'X + (1/h)A]^{-1}X'Y$$

Proof. Since $A$ is p.d. we can factor $A$ as $A = S'S$. Condition (2) becomes $\gamma'\gamma \leq \sigma^2 h$ where $S\beta = \gamma$. Hence, $\hat{\sigma}$ must minimize

$$\max_{\gamma'\gamma \leq \sigma^2 h} \quad E_{\gamma,\sigma}[(\sigma'Y - a'S^{-1}\gamma)^2]$$

$$= \sigma^2 \sigma'\sigma + \max_{\gamma'\gamma \leq \sigma^2 h}[(\sigma'XS^{-1}\gamma - a'S^{-1}\gamma)^2]$$

$$= \sigma^2 \sigma'\sigma + \sigma^2 h \max_{\gamma'\gamma \leq 1}[((\sigma'XS^{-1} - a'S^{-1})\gamma)^2]$$

$$= \sigma^2 \sigma'\sigma + \sigma^2 h (S^{-1'}X'\sigma - a')A^{-1}(X'\sigma - a) \qquad (5)$$

-4-

by Lemma 1. By finding the zero of the vector derivative of (3) with respect to $\sigma$, one can see that the $\sigma$ which minimizes (3) is given by

$$
\begin{aligned}
\hat{\sigma} &= [\dots + XA^{-1}X'J^{-1}XA^{-1}\alpha \\
&= X[\dots + A^{-1}X'XJ^{-1}A^{-1}\alpha \\
&= \dots + (\dots)A J^{-1}\alpha
\end{aligned}
$$

(since $\dots[\dots + A^{-1}X'XJ = [(1/\lambda)I_N + XA^{-1}X'JX)$
where the subscript of $I$ denotes the order of the identity matrix.
The theorem follows immediately, since the Hessian of (3) is always positive definite.

It follows from Theorem 1, by letting $c$ be the Kronecker delta $\delta_{ij}$, that minimax estimation of individual components $\sigma_i$ of $\sigma$ based on (2) also produces the estimator $\hat{\sigma}_{m,i}$ of $\sigma$.

One can easily verify that the minimax linear estimator $\hat{\sigma}_{m,i}$ based on (2) is equivalent to the minimax linear estimator based on

$$
\beta'A\beta/\sigma^2 = \lambda. \tag{4}
$$

Hence, minimax theory suggests that if $\hat{\lambda}$ is an estimator of $\beta'A\beta/\sigma^2$, then one might use

$$
\hat{\sigma}_{m,\hat{\lambda}} = [A'A + (1/\hat{\lambda})A]^{-1}X'Y
$$

to estimate $\Delta$. One could also argue that overestimation of $\Delta$ by $\hat{\Delta}$ seems preferable to underestimation since, in (2), $\Delta$ is an upper bound; however, robustness results discussed in [11] and simulation results given in the present paper indicate that overestimation may result in unnecessarily and undesirably cautious estimators.

Example 1. Consider the case when $A = I_q$. The minimax linear estimator $\hat{\beta}_{L,\Delta}$ is identical to the ridge regression estimator $\hat{\beta}_{k}$ where $k = \Delta^{-1}$. Minimax theory supports the estimator $\hat{\beta}_{L,\hat{\Delta}}$ where $\hat{\Delta} = \hat{\beta}'\hat{\beta}/\hat{\sigma}^2$ and $\hat{\beta}$ and $\hat{\sigma}^2$ are the usual estimators of $\beta$ and $\sigma^2$ given by $\hat{\beta} = (X'X)^{-1}X'Y$ and $\hat{\sigma}^2 = (Y-X\hat{\beta})'(Y-X\hat{\beta})/n-q$.

Example 2. Let $A = X'X$. Since $A$ is p.d., we can write $A = S'S$ where $S$ is a nonsingular $q \times q$ matrix. Hence, the model $Y = X\beta + \varepsilon$ can be written as $Y = \tilde{X}\tilde{\beta} + \varepsilon$ where $\tilde{X} = XS^{-1}$ and $\tilde{\beta} = S\beta$. Since $\tilde{X}'\tilde{X} = I$, it seems reasonable to place equal weight on the size of different components of $\tilde{\beta}$; that is, the constraint $\tilde{\beta}'\tilde{\beta} \leq \Delta$ seems very reasonable. This constraint is equivalent to the constraint $\beta'X'X\beta \leq \Delta$ which leads to $\hat{\beta}_{L,\Delta} = [X'X + (1/\Delta)X'X]^{-1}X'Y = (1 + \Delta^{-1})^{-1}\hat{\beta}$. This estimator is a shrunken estimator. Estimation of $\Delta$ by $\hat{\Delta} = \hat{\beta}'X'X\hat{\beta}/\hat{\sigma}^2$ leads to the estimator $(1 + \hat{\Delta}^{-1})^{-1}\hat{\beta}$ which is a shrunken estimator with a stochastic shrinkage coefficient.

Example 3. One also might also consider minimax linear regression estimators based on constraints of the form

$$(\beta - \beta^*)'A(\beta - \beta^*), \sigma^2 \leq \Delta$$

-13-

where $A$, $\alpha$, and $h$ are known. Some of the results from [9] and some additional related results are given in [2]. See Chapter 10 of [11] for a detailed, translated discussion of the results of [9].

Kuks and Olman show in [9] that for any $a \in R^q$,

$$E_{b,\sigma}[(a'\hat{\beta}_{A,\alpha} - a'\beta)^2] \le E_{\beta,\sigma}[(a'\hat{\beta} - a'\beta)^2]$$

for all $\beta, \sigma$ satisfying $\beta'A\beta/\sigma^2 \le 2h$. A generalization in [11] of a result from [9] states that for any nonnegative definite (n.d.) $q \times q$ matrix $W$,

$$E_{b,\sigma}[(\hat{\beta}_{A,\alpha} - \beta)'W(\hat{\beta}_{A,\alpha} - \beta)] \le E_{\beta,\sigma}[(\hat{\beta} - \beta)'W(\hat{\beta} - \beta)]$$

for all $\beta, \sigma$ satisfying $\beta'A\beta/\sigma^2 \le 2h$. Peele observes in [11] that the two above-mentioned results, which imply robustness of $\hat{\beta}_{A,\alpha}$ relative to assumption (2), are equivalent and notes that the $A$ matrices in Examples 1 and 2 are optimal (in a most-robust sense) prior information matrices. Clearly the accuracy of $\hat{\beta}_{A,\alpha}$ as an estimator of $\beta$ will depend on the accuracy of the prior information (2) which involves $A$ and $h$.

The following invariance property of minimax linear estimation supports the model-transformation argument used in Example 2 and also answers a criticism of (non-minimax) ridge regression techniques made in [11].

-7-

*Theorem 1.* Let $T$ be any $q \times q$ nonsingular matrix. It follows that the model $Y = X\beta + \xi$ is equivalent to the model $Y = \tilde{X}\tilde{\beta} + \xi$ where $\tilde{X} = XT$ and $\tilde{\beta} = T^{-1}\beta$. Also, the condition $\beta' A\beta/\sigma^2 \leq k$ is equivalent to the condition $\tilde{\beta}' T' A T\tilde{\beta}/\sigma^2 \leq k$ since $\beta = T\tilde{\beta}$. Minimax linear estimation is invariant under model transformations in that $\hat{\beta}_{A,k} = T\hat{\tilde{\beta}}_{A,k}$.

*Proof.* One can easily check that

$$
\begin{aligned}
\hat{\tilde{\beta}}_{A,k} &= [\tilde{X}'\tilde{X} + (1/k)T'AT]^{-1}\tilde{X}'Y \\
&= [T'X'XT + (1/k)T'AT]^{-1}T'X'Y \\
&= [X'X + (1/k)A]^{-1}X'Y \\
&= \hat{\beta}_{A,k} .
\end{aligned}
$$

The following theorem appears to answer a question posed in Chapter II of [1] concerning the existence of a characterization of minimax linear regression estimation similar to the least squares characterization of BVUE regression estimators.

*Theorem 2.* Among all $b \in \mathbb{R}^q$, $b = \hat{\beta}_{A,k}$ minimizes

$$
(1/k)b'Ab + (Y - Xb)'(Y - Xb) . \tag{5}
$$

*Proof.* The zero of the vector derivative of (5) with respect to $b$ is $b = \hat{\beta}_{A,k}$, and the Hessian of (5) is always positive definite.

-8-

An immediate consequence of Theorem 3 is the result that among $b$ satisfying $b'Ab = (\hat{\beta}_{A,k})'A(\hat{\beta}_{A,k})$, $(Y - Xb)'(Y - Xb)$ is minimized by $b = \hat{\beta}_{A,k}$. This result has been previously observed in [7] for the case $A = I$ by Lagrange multiplier methods. See [16] for a discussion of results similar to Theorem 3.

One can easily check by orthogonal transformation methods similar to those on page 62 of [7] that if

$$g_1(k) = [\hat{\beta}^*(k)]'[\hat{\beta}^*(k)]$$

then $g_1$ is a decreasing function of $k$. Hence, it follows easily (by supposing not and then contradicting Theorem 3) that if

$$g_2(k) = [Y - X\hat{\beta}^*(k)]'[Y - X\hat{\beta}^*(k)]$$

then $g_2$ is an increasing function of $k$.

Two loss functions different from $g_2(k)$, which is minimized by the least squared estimator $\hat{\beta} = \hat{\beta}^*(0)$, are considered in Section 3 and simulation results are given for minimax-related ridge regression estimators. The simulations involve progressive ill-conditioning and some of the simulations use random, norm-one $\beta$.

## 3. NUMERICAL RESULTS

Within the past five years, many ridge estimators have appeared in the literature and have been compared in subsequent simulation studies. Golub, Heath, and Wahba [5] state that two dozen is probably a conservative estimate of the number of such estimators.

Many practitioners have used the ridge-trace approach to select the value of the ridge constant. This approach is highly subjective, however, and has also been criticized on other grounds (see, e.g., Smith and Campbell [14]). Consequently, it seems as though more formal methods for selecting the ridge constant need to be employed. The ridge estimator presented in Section 2,

$$\hat{\beta}*(\hat{k}_1) = (X'X + \hat{k}_1 I)^{-1} X'Y$$

where $\hat{k}_1$ is stochastic and equal to $\hat{\sigma}^2 / \hat{\beta}'\hat{\beta}$, has been the focal point in our simulation studies. Some of these results are presented later in this section. We have also studied the shrinkage estimator presented in Example 2, but it did not fare as well as the ridge estimator.

Before undertaking a simulation study of regression estimators, an experimenter must select an appropriate loss function, and choose a mechanism for generating data that is representative of actual data. The loss function most frequently used is squared error, i.e.

$$(\tilde{\beta}-\beta)'(\tilde{\beta}-\beta) \hspace{3cm} (6)$$

where $\tilde{\beta}$ denotes an estimator of $\beta$ (see, e.g. Gibbons [4] and Hemmerle and Brantle [6]).

We have used (6) in addition to the loss function

$$(\tilde{\beta}-\beta)'X'X(\tilde{\beta}-\beta), \hspace{3cm} (7)$$

which was also used by Dempster, Schatzoff, and Wermuth [3]. It can be argued (as they did) that (6) would be an appropriate loss function if the primary goal of a regression study is estimation of the parameters, whereas (7) would be appropriate if the regression was to be used for prediction (since (7) may be written as $(X\tilde{\beta}-X\beta)'(X\tilde{\beta}-X\beta)$, and $Y = X\beta + \varepsilon$ is predicted by $\hat{Y} = X\tilde{\beta}$). The loss function $(Y-X\tilde{\beta})'(Y-X\tilde{\beta})$ would be inappropriate for either estimation or prediction since the loss is minimized when $\tilde{\beta}$ is the least squares estimator, as was mentioned in Section 2.

If we adopt either (6) or (7) or both, we then must decide what $X$ and $\beta$ to use as well as the method for generating values of $Y$. Newhouse and Oman [10] showed that, assuming $X$, $\sigma^2$, and $k$ to be fixed, $MSE(\hat{\beta}^*(k))$ is maximized, for $\beta'\beta = 1$, when $\beta$ is the normalized eigenvector corresponding to the smallest eigenvalue of $X'X$. Similarly, $MSE(\hat{\beta}^*(k))$ is minimized when $\beta$ is the normalized eigenvector corresponding to the largest eigenvalue of $X'X$. This result has apparently led researchers to

-11-

use these two choices for $\beta$ in their simulation studies. Typically, $X$ is selected in such a way as to make $X'X$ an equicorrelation matrix, and observations on the dependent variable are then generated as

$$Y = X\beta + e \qquad (8)$$

where $e$ is $N(0, \sigma^2)$. Several values of $\sigma^2$ are then used in conjunction with each of several $X$ matrices. The loss for each estimator is then determined using either (6) or (7), or both.

This general procedure has two major shortcomings. First, we would not expect to encounter an equicorrelation matrix with actual data. Second, although the two choices for $\beta$ correspond to the best and to the worst case for ridge regression estimators (in terms of mean squared error), we advise against using the smallest eigenvalue case. The reason can be discerned from Table 1. In particular, we should notice what happens to the average value of $R^2$ (the coefficient of multiple determination) as we move, keeping $\sigma^2$ fixed, from the well-conditioned matrix in the first part of the table, to the highly ill-conditioned matrix in the bottom part of the table. One can observe (as in Table 1) that $R^2$ will be approximately $q/n$ in the smallest eigenvalue case even for moderate $\sigma^2$ if $X'X$ is highly ill-conditioned. Since this is the worst case for ridge estimators,

-12-

we are thus unable to see how poorly a ridge estimator performs relative to least squares for reasonable values of $R^2$. Since progressive ill-conditioning tends to cause $R^2$ for the smallest eigenvalue case to be much different from $R^2$ for the largest eigenvalue case, making $\sigma^2$ smaller so as to produce higher $R^2$ values for the smallest eigenvalue case would tend to make $R^2$ almost exactly 1.0 in the largest eigenvalue case.

For these reasons, we have used a different simulation procedure. We have also generated observations on $Y$ as in (8), but, for each trial, a $\beta$ vector is generated from a uniform distribution on the collection of all norm-one q-element vectors.

A normal$(0, \sigma^2)$ error vector $e$ is then generated and $Y$ is computed as $Y = X\beta + e$. Thisted [10] cautions against the use of random $\beta$ in simulation studies which include different classes of estimators. We see nothing wrong, however, with using $\beta$ vectors that are uniform on the collection of all norm-one vectors when comparing estimators within a particular class such as ridge estimators. As will be seen later, the main objectives of our numerical studies have been: (1) to compare our ridge estimator with the much-referenced Hoerl, Kennard, and Baldwin [8] estimator, and (2) to determine under what conditions our ridge constant or a multiple of our ridge constant would be appropriate.

In addition to our different procedure for generating $\beta$,

-13-

we have also generated $X'X$ matrices in a different manner. Instead of the equicorrelation approach, we used the method given in Ryan [13] to generate several progressively ill-conditioned matrices.

The results of our major simulation study are shown in Table 2. The following can be discerned from inspection of the table. For a particular degree of ill-conditioning, the size of the desired ridge constant depends on the size of $R^2$ in that a large $R^2$ implies the need for a small ridge constant, and vice versa. Secondly, for fixed $R^2$, the size of the desired ridge constant varies with the degree of ill-conditioning in that greater ill-conditioning implies the need for a larger ridge constant.

It follows from the discussion in Section 2 concerning robustness that the use of a nonstochastic ridge constant $k$ will result in the ridge estimator outperforming least squares in terms of MSE provided that $k < (h/2)^{-1} = 2h^{-1}$ where $h = \beta'\beta/\sigma^2$. Consequently, $2\hat{h}^{-1} = 2\hat{\sigma}^2/\hat{\beta}'\hat{\beta}$ (in addition to $\hat{h}^{-1}$ given in Example 1) should be a good (stochastic) choice for $k$. This choice can be motivated, in part, by the results in Table 2, and also by the fact that $E(\hat{\sigma}^2/\hat{\beta}'\hat{\beta})$ should be less than $E(\sigma^2/\beta'\beta)$ since (assuming normality) $E(\hat{\beta}'\hat{\beta}) = \beta'\beta + \sigma^2\sum_{i=1}^{i}\lambda_i^{-1}$, where $\lambda_1, \lambda_2, \ldots, \lambda_q$ are the eigenvalues of $X'X$ and $\hat{\sigma}^2$ is an unbiased estimator of $\sigma^2$.

One might argue that $\hat{\beta}'\hat{\beta}$ in $\hat{h}_1$ should be replaced by $\hat{\beta}'\hat{\beta} - \hat{\sigma}^2\sum_{i=1}^{}\lambda_i^{-1}$ so as to make $\hat{h}_1$ approximately unbiased. When $X'X$ is in correlation

form $\sigma^2$ must be quite small to produce good $R^2$ values. Accordingly, $\hat{\lambda}_1^{-1}$ would remain essentially unchanged unless $X'X$ is highly ill-conditioned. In particular, the last matrix in Table 2 is just such a matrix. For that matrix, with $\hat{\lambda}_2 = (\hat{\beta}'\hat{\beta} - \sigma^2_{jj}\lambda_j^{-1})/\hat{\beta}^2$, the average value of $\sigma\hat{\lambda}_2^{-1}$ should be approximately $\sigma\hat{\lambda}_1^{-1}$ for $\sigma^2 = .0001$, approximately $\sigma\hat{\lambda}_1^{-1}$ for $\sigma^2 = .001$, and approximately $0.8\,\hat{\lambda}_1^{-1}$ for $\sigma^2 = .005$. Notice that these are close to the minimum values in Table 2, especially for the second loss function. Therefore, especially if $X'X$ is highly ill-conditioned, $\hat{\lambda}_1^{-1}$ or $(\hat{\lambda}_2/\sigma)^{-1}$ be considered if the numerator of $\hat{\lambda}_2$ is sufficiently positive. Although the vast majority of our simulations have been for $q = 4$, limited work with $q = 3$ has revealed that $\sigma\hat{\lambda}_2^{-1}$ may be a good choice for $k$ when $X'X$ is highly ill-conditioned and $R^2$ is very close to 1.0.

Many researchers have contended that the improvement of a biased estimator over least squares will usually be quite small. The figures in Table 2 indicate that this is not necessarily true. The absolute gain might be small since the data for the regressors is in correlation form, but could be quite sizable after converting back to raw form.

TABLE 1 — Comparison of least squares, the minimum estimator ($k = \hat{h}_1^{-1}$), and the Hoerl, Kennard, and Baldwin estimator ($k = 4\hat{h}_1^{-1}$).

Well-conditioned X matrix ($\lambda_- = .7796$, $\lambda_+ = 1.4417$)

| $\sigma^2$ | Ave. $R^2$ | Least Squares | $\hat{h}_1^{-1}$ | $4\hat{h}_1^{-1}$ |
|---|---|---|---|---|
| .0001 | .9954 | .000410 | .000410 | .000410 |
| | .9976 | .000396 | .000396 | .000395 |
| .001 | .9564 | .005008 | .004993 | .004968 |
| | .9753 | .004349 | .004334 | .004295 |
| .01 | .5949 | .044814 | .044376 | .044519 |
| | .5014 | .042076 | .041355 | .039948 |
| .1 | .2441 | .465888 | .434394 | .401581 |
| | .3201 | .392342 | .346304 | .286141 |

Moderately ill-conditioned X matrix ($\lambda_- = .1250$, $\lambda_+ = 3.2700$)

| $\sigma^2$ | Ave. $R^2$ | Least Squares | $\hat{h}_1^{-1}$ | $4\hat{h}_1^{-1}$ |
|---|---|---|---|---|
| .0001 | .9721 | .001655 | .001657 | .001675 |
| | .9984 | .001852 | .001850 | .001872 |
| .001 | .9721 | .014137 | .014990 | .017929 |
| | .9976 | .017555 | .017630 | .016965 |
| .01 | .9521 | .155802 | .151254 | .173569 |
| | .9885 | .153610 | .145073 | .111957 |
| .1 | .1976 | 1.673.21 | 1.28927 | .953034 |
| | .9963 | 1.73359 | 1.21755 | .617053 |

Very ill-conditioned X matrix ($\lambda_- = .0031$, $\lambda_+ = 3.6372$)

| $\sigma^2$ | Ave. $R^2$ | Least Squares | $\hat{h}_1^{-1}$ | $4\hat{h}_1^{-1}$ |
|---|---|---|---|---|
| .0001 | .4797 | .029231 | .032594 | .030302 |
| | .9990 | .029267 | .027737 | .023895 |
| .001 | .1628 | .379005 | .428802 | .498001 |
| | .9902 | .374973 | .278236 | .151009 |
| .01 | .0963 | 2.89948 | 1.51099 | 1.01671 |
| | .9059 | 2.84197 | 1.35998 | .488576 |
| .1 | .0960 | 28.8260 | 13.5135 | 5.66249 |
| | .5179 | 28.2283 | 12.8263 | 4.91546 |

The loss is $(\hat{\beta}_{est} - \beta)'(\hat{\beta}_{est} - \beta)$. The X matrix is $40 \times 4$; $\lambda_+$ and $\lambda_-$ denote the largest and smallest eigenvalues, respectively, of X'X. The first row of averages is for the smallest eigenvalue case; the second row for the largest. Averages are over 100 observations.

TABLE 2 — Comparison of minimax-related ridge estimators using random norm--ing B.

## Well-conditioned X matrix

| $\sigma^2$ | Avg. $R^2$ | Least Squares | $\hat{h}_1^{-1}$ | $2\hat{h}_1^{-1}$ | $4\hat{h}_1^{-1}$ | $6\hat{h}_1^{-1}$ | $8\hat{h}_1^{-1}$ | $10\hat{h}_1^{-1}$ | Loss |
|---|---|---|---|---|---|---|---|---|---|
| .0001 | .996 | .0004218 | 1.0002 | 1.0007 | 1.0017 | 1.0031 | 1.0045 | 1.0062 | A |
|  |  | .0004020 | 1.0005 | 1.0007 | 1.0017 | 1.0030 | 1.0042 | 1.0060 | B |
| .001 | .964 | .004425 | 0.9991 | 0.9984 | 0.9991 | 1.0020 | 1.0070 | 1.0142 | A |
|  |  | .004140 | 0.9988 | 0.9983 | 0.9986 | 1.0015 | 1.0056 | 1.0123 | B |
| .005 | .846 | .02060 | 0.9961 | 0.9951 | 1.0015 | 1.0154 | 1.0456 | 1.0811 | A |
|  |  | .01939 | 0.9959 | 0.9948 | 1.0008 | 1.0160 | 1.0407 | 1.0743 | B |

## Moderately Ill-conditioned X matrix

| $\sigma^2$ | Avg. $R^2$ | Least Squares | $\hat{h}_1^{-1}$ | $2\hat{h}_1^{-1}$ | $4\hat{h}_1^{-1}$ | $6\hat{h}_1^{-1}$ | $8\hat{h}_1^{-1}$ | $10\hat{h}_1^{-1}$ | Loss |
|---|---|---|---|---|---|---|---|---|---|
| .0001 | .993 | .001831 | 1.0001 | 1.0005 | 1.0024 | 1.0055 | 1.0099 | 1.0156 | A |
|  |  | .0004198 | 1.0002 | 1.0007 | 1.0021 | 1.0045 | 1.0076 | 1.0117 | B |
| .001 | .937 | .01842 | 0.9962 | 0.9957 | 1.0038 | 1.0233 | 1.0527 | 1.0917 | A |
|  |  | .004233 | 0.9979 | 0.9979 | 1.0038 | 1.0172 | 1.0377 | 1.0644 | B |
| .005 | .854 | .0537 | 0.962 | 0.935 | 0.912 | 0.913 | 0.939 | 0.977 | A |
|  |  | .0196 | 0.974 | 0.959 | 0.944 | 0.954 | 0.974 | 1.005 | B |

## Very Ill-conditioned X matrix

| $\sigma^2$ | Avg. $R^2$ | Least Squares | $\hat{h}_1^{-1}$ | $2\hat{h}_1^{-1}$ | $4\hat{h}_1^{-1}$ | $6\hat{h}_1^{-1}$ | $8\hat{h}_1^{-1}$ | $10\hat{h}_1^{-1}$ | Loss |
|---|---|---|---|---|---|---|---|---|---|
| .0001 | .987 | .0326 | 0.982 | 0.973 | 1.000 | 1.053 | 1.129 | 1.221 | A |
|  |  | .000380 | 0.997 | 0.995 | 1.005 | 1.024 | 1.050 | 1.079 | B |
| .001 | .932 | .339 | 0.796 | 0.630 | 0.590 | 0.555 | 0.543 | 0.546 | A |
|  |  | .00398 | 0.942 | 0.912 | 0.867 | 0.899 | 0.920 | 1.000 | B |
| .005 | .862 | 1.631 | 0.604 | 0.452 | 0.318 | 0.259 | 0.228 | 0.210 | A |
|  |  | .0199 | 0.874 | 0.824 | 0.769 | 0.749 | 0.739 | 0.990 | B |

600 "good" ($R^2 > .7$) observations were generated for each $\sigma^2$, matrix combination. The entries represent average loss for least squares and the ratio of average loss for each estimator to the average loss for least squares.

Loss A is $(\hat{\beta}_{RR} - \beta)'(\hat{\beta}_{RR} - \beta)$        Loss B is $(\hat{\beta}_{RR} - \beta)'X'X(\hat{\beta}_{RR} - \beta)$

$$= (\hat{Y} - X\beta)'(\hat{Y} - X\beta)$$

The matrices used here are the same as those used in TABLE 1.

## 4. SUMMARY

A minimax approach to linear regression was presented and some of the pertinent minimax results of an important Russian paper by Kuks and Olman were also discussed. Additional minimax results were presented including an invariance property of minimax linear regression estimators. In Section 3 some criticisms were made of the most common simulation procedure for ridge regression comparisons, and these were also illustrated in Table 1. Accordingly, a different procedure involving random $\beta$ was used to produce Table 2.

The minimax regression estimators with ridge constants $\hat{\lambda}^{-1}$ and $s\hat{\lambda}^{-1}$ were seen to perform well in the major simulation study although $s\hat{\lambda}^{-1}$ may be a better choice, especially if $X'X$ is highly ill-conditioned.

## 5. ACKNOWLEDGMENTS

REFERENCES

[1] BIBBY, J. and TOUTENBURG, H. (1977). Prediction and
    Improved Estimation in Linear Models. New York: Wiley.

[2] BUNKE, O. (197 ). Improved inference in linear models with
    additional information. Math. Op. and Stat., , 817-829.

[3] DEMPSTER, A., SCHATZOFF, M., and WERMUTH, N. (1977). A
    simulation study of alternatives to ordinary least squares.
    J. Amer. Statist. Assoc., 72, 77-90.

[4] GIBBONS, D. (197 ). A simulation study of some ridge
    estimators. General Motors Research Laboratories, Research
    Publication GMR-2659. (Revised), September.

[5] GOLUB, G., HEATH, M., and WAHBA, G. (1977). Generalized
    cross-validation as a method for choosing a good ridge
    parameter. Technometrics, 21, 215-223.

[6] HEMMERLE, W. J. and BRANTLE, T. F. (197 ). Explicit and
    constrained generalized ridge estimation. Technometrics, ,
    109-120.

[7] HOERL, A. E., and KENNARD, R. W. (197 ). Ridge regression:
    biased estimation for nonorthogonal problems. Technometrics,
    12, 55-67.

[8] HOERL, A. E., KENNARD, R. W., and BALDWIN, K. F. (1975).
    Ridge regression: some simulations. Comm. in Statist., 4,
    105-12 .

[9]  KUKS, J. and OLMAN, W.  (1972).  Minimax linear estimation
     of regression coefficients, II.  Izvestija Akademii Nauk
     Estonskoj SSR 1., 66-72.

[10] Newhouse, J. P. and OMAN, S. D.  (1971).  An evaluation of
     ridge estimators.  Rand Corporation Report No. R-716-PR.

[11] PEELE, L.  (1979).  Most-robust ridge regression estimators.
     Old Dominion University Department of Mathematical Sciences
     Technical Report TR 79-7, September.

[12] ROYDEN, H.  (1968).  *Real Analysis*.  New York:  Macmillian.

[13] RYAN, T. P.  (1980).  A new method of generating correlation
     matrices.  *J. Statist. Comput. Simul.*, to appear.

[14] SMITH, G. and CAMPBELL, F.  (1980).  A critique of ridge
     regression methods.  *J. Amer. Statist. Assoc.*, to appear.

[15] THEOBALD, C. M.  (1974).  Generalizations of mean squared
     error applied to ridge regression.  *J. Roy. Statist. Soc.*,
     B, 36, 103-106.

nature. Traces of MSE estimates and of related direction cosine estimates can be used to augment the trace of ridge coefficients. It is the "see power" of these data analytic displays that glorifies enlistment in the ridge navy.

## 2. INVARIANCE

Now consider the effect on variance of a linear model as indicated by the expression of a vector invariant $X\beta = Z\gamma$ for $Z = XA$ and $\gamma = A^{-1}\beta$. But the amount of ridge shrinkage that is MSE optimal in the usual trace sense (Obenchain 1975) is a non-linear function of the eigenvalues of $X'X$ (or $Z'Z$) and of the parameters, $\beta$ (or $\gamma$) and $\sigma^2$. The unknown MSE optimal ridge fit is thus not invariant. Why, then, do Smith and Campbell feel that the fit "should be" invariant?

An estimation procedure will yield fits that are invariant in this sense if and only if the estimate $b$ of $\beta$ and the estimate $c$ of $\gamma$ are interrelated by the implicit formula $c = A^{-1}b$ for every nonsingular $A$. The ridge solutions that optimize the objective criteria of my earlier comments are not implicitly interrelated. And I feel that no "anxiety" worth avoiding is implied by this property.

It is well known that any procedure that yields implicitly interrelated estimates cannot have certain desirable properties. For example, Tukey (1975) considers the two-regressor case $z_1 = x_1 + ax_2$ and $z_2 = x_2$. The resulting implicit estimates are $c_1 = b_1$ and $c_2 = b_2 - ab_1$. This represents a change in the implicit estimate of the unchanged regressor and no change in the implicit estimate of the changed regressor. In this sense, no implicit coefficient can be interpreted as measuring the "effect" of only one regressor on the expected response.

Ridge estimates can be chosen to be less highly intercorrelated than are the implicit least squares estimates for an ill-conditioned parameterization. These ridge estimates are thus "a step in the right direction" in the sense of Tukey (1975) toward coefficients that can be interpreted as "effects" in the sense of Obenchain and Vinod (1974).

## 3. OBTAINING PREASSIGNED VALUES

Smith and Campbell make a very strong point when they assert that, if ridge regression is assumed to yield implicit estimates, then these coefficients can be made to correspond to any preassigned vector by choice of a diagonal matrix $A$ in their equation (2.2).

If one considers generalized ridge regression with arbitrary origin, as in Smith and Campbell's equation (5.1), then one can indeed obtain any preassigned solution by choice of the ridge factors, $a_1, \ldots, a_n$. But, instead of $0 \le a_i \le 1$, one might have to go to the extreme of using ridge factors greater than one or even negative. If one is monitoring sample information on the variance-bias trade-off, then one would see that MSE is not reduced by such an extreme tactic.

## REFERENCES

Efron, B., and Morris, C. (1977), "Comment," *Journal of the American Statistical Association*, 72, 91-93.

Mallows, C.L. (1973), "Some Comments on $C_p$," *Technometrics*, 15, 661-675.

McCabe, G.P. (1978), "The Selection of Regression Coefficients Using Acceptability," *Technometrics*, 20, 131-139.

Obenchain, R.L. (1975), "Ridge Analysis Following a Preliminary Test of the Shrunken Hypothesis," *Technometrics*, 17, 431-441 (with discussion by G.C. McDonald, 443-44).

——— (1977), "Classical F-tests and Confidence Regions for Ridge Regression," *Technometrics*, 19, 429-439.

——— (1978), "Good and Optimal Ridge Estimators," *Annals of Statistics*, 6, 1111-1121.

Obenchain, R.L., and Vinod, H.D. (1974), "Estimates of Partial Derivatives From Ridge Regression," in *Invited paper at NBER-NSF Seminar on Bayesian Inference in Econometrics*, Ann Arbor, Mich.

Tukey, J.W. (1975), "Instead of Gauss-Markov Least Squares, What?," in *Applied Statistics*, ed. R.P. Gupta, Amsterdam: North-Holland Publishing Co., 352-372.

# Comment

## LAWRENCE C. PEELE and THOMAS P. RYAN*

The authors have contributed some interesting criticisms but their contentions seem not to form a basis for general condemnation of ridge regression. Smith and Campbell state that they advocate the use of prior information, but are uncomfortable with exact linear constraints. They also state that "when the least squares estimates are imprecise, auxiliary information is quite useful. Pseudoinformation is of dubious value."

* Lawrence C. Peele is Assistant Professor in the Department of Mathematical Sciences, Old Dominion University, Norfolk, VA 23508. Thomas P. Ryan is Assistant Professor in the Department of Quantitative Analysis, University of Cincinnati, Cincinnati, OH 45221. Research was supported in part by the Air Force Office of Scientific Research, Contract No. AFOSR F49620-79-C-0125.

We agree that exact linear constraints are "discomforting" (and seemingly unrealistic, as well) and, consequently, prefer constraints of the form $\beta'\beta \leq h$, or more generally $\beta'T\beta \leq h\sigma^2$, where $T$ is a positive definite matrix and $h$ is a positive constant. There is no boundedness assumption on $\beta'\beta$ inherent in the least squares procedure, and, in practice, one might set an a priori upper bound on reasonable values of $\beta'\beta$. The degree of superiority of a biased estimator over least squares will naturally depend on how sharp a bound on $\beta'\beta$ an experimenter can impose. Suppose that based on prior information only, one believes that $\sigma^{-2}\beta'\beta$ is approximately equal to a number $h^*$. Kuks and Olman (1972) showed that if $h^* \geq \sigma^{-2}\beta'\beta/2$, then the ridge regression estimator with ridge parameter $k = 1/h^*$ will outperform the least squares estimator in terms of mean squared error. Consequently, if fairly accurate prior information about $\sigma^{-2}\beta'\beta$ is available, this ridge estimator is preferable to the least squares estimator. Lacking prior information, we can still obtain, by estimating $\sigma^{-2}\beta'\beta$ from sample data, an estimator that tends to be better than least squares. It follows from this line of reasoning that estimators with this underlying motivation will not have "loose theoretical underpinnings."

Smith and Campbell note that ridge regression estimators are not invariant under model transformations of the form

$$Y = (XA)A^{-1}\beta + \epsilon = Z\gamma + \epsilon. \qquad (1)$$

Peele and Ryan (1979) discuss minimax linear estimation based on prior information of the form $\beta'T\beta \leq \sigma^2 h$, which includes ordinary and generalized ridge regression estimation. The authors observe that the resulting estimators

$$\hat{\beta}^*_{T,h} = \left[\frac{1}{h} T + X'X\right]^{-1} X'Y$$

are invariant under model transformations of the form (1).

In his contribution to this discussion, Van Nostrand briefly mentioned specific (conflicting) results from several simulation studies. It is not surprising that the results can differ greatly from one study to the next since no "standard" simulation procedure is being used. In point of fact, it seems that those of us who have performed these studies have fallen to a number of pitfalls. In a regression simulation study, one should use appropriate values of the error variance $\sigma^2$. The use of correlation form implies that one should choose $\sigma^2$ close to zero in order to generate data with reasonable $R^2$ values. Unfortunately, $R^2$ is quite sensitive to small changes in $\sigma^2$; consequently, great care must be exercised in order to avoid producing small $R^2$ values. Also, the usual procedure of letting $\beta$ be represented by the eigenvectors corresponding to the largest and smallest eigenvalues of $X'X$ is objectionable. In the smallest eigenvalue case, $R^2$ will be approximately $p/n$ when the smallest eigenvalue is quite small and $\sigma^2$ is at least one order of magnitude larger. Thus, with $p = 4$ and $n = 100$ (a typical choice) we could hardly say that we have generated representative regression data when $R^2$ is approximately .04. Making $\sigma^2$ smaller in an effort to remedy the problem would only cause $R^2$ to be almost exactly 1.0 in the largest eigenvalue case. There is clearly a need for a better approach.

## REFERENCES

Kuks, J., and Olman, W. (1972), "Minimax Linear Estimation of Regression Coefficients, II," *Izvestiya Akademiya Nauk Estonskoy SSR* 21, 66-72.

Peele, L., and Ryan, T. (1979), "Minimax Regression Estimators With Application to Ridge Regression," submitted to *Technometrics*.

# Comment

## H. D. VINOD*

Certain prior knowledge about intrinsic measurement errors can be cleverly incorporated in ridge regression by adding $p$ fictitious observations to the data. Consider a standardized regression model $y = X\beta + u$, where $X'X$ is the correlation matrix based on $n$ observations among $p$ regressors and $X'y$ is a vector of correlation coefficients with the dependent variable. There are

measurement errors in each regressor, so that the available data are indistinguishable from other data where one may add a number between $-.5$ to $.499$ beyond the last published digit. There may be a similar error in each mean ($\bar{x}_i$) and standard deviation (SD). The largest measurement error ($e^{max}$) in each standardized regressor $(x_{it} - \bar{x}_i)(n - 1)^{-1/2}(SD_i)^{-1}$, where $it = 1, 2, \ldots, n$, is

* H. D. Vinod is Supervisor, Economic Studies, Economic Analysis Section, American Telephone and Telegraph Co., Piscataway, NJ 08854.

# DATE
# FILMED

8